

الگوریتم بازیابی و رتبه‌بندی اطلاعات در موتور جستجوی گوگل

سعیده ابراهیمی^۱

چکیده

هدف اساسی این مقاله، تبیین نحوه عمل موتور جستجوی گوگل در رتبه‌بندی اطلاعات بازیابی شده است و بدین منظور، الگوریتم (خوارزمی) موتور جستجوی گوگل را در بازیابی و رتبه‌بندی اطلاعات بررسی می‌کند. در بخش اول - که text matching نامیده می‌شود - شیوه یافتن اطلاعات مرتبط با واژه‌های وارد شده در جعبه جستجو، به‌طور خلاصه بیان می‌گردد و در بخش دوم - که بخش اصلی مقاله است و PageRank نام دارد - شیوه رتبه‌بندی نتایج مرتبط بازیابی می‌شود و به تفصیل مورد بررسی قرار می‌گیرد. در ادامه، الگوی موج‌سوار تصادفی، به‌عنوان تأییدی بر الگوریتم (خوارزمی) رتبه‌بندی گوگل تشریح می‌گردد. این مقاله در پایان، روش مشخص کردن رتبه صفحات وب را از طریق نوار ابزار گوگل توصیف می‌کند.

کلیدواژه‌ها

موتور جستجوی گوگل، الگوریتم، PageRank، رتبه‌بندی اطلاعات، لینک، الگوی موج‌سوار تصادفی.

مقدمه

جستجو در آن نیز تا اندازه زیادی بر این مسئله سایه انداخته است. گوگل، هم برای کاربران و هم برای مدیران وبگاه‌ها دارای ارزشی منحصر به فرد است. خزنده‌های موجود در موتور گوگل،

در طی چند سال اخیر، گوگل به پراستفاده‌ترین موتور جستجوی جهان تبدیل شده است. این مسئله، نه تنها به عملکرد و سهولت استفاده از این موتور برمی‌گردد، بلکه کیفیت نتایج

۱. دانشجوی دکتری کتابداری و اطلاع‌رسانی دانشگاه شهید چمران اهواز. ebrahimi_saedeh@yahoo.com

با سرعت قابل ملاحظه‌ای، وبگاه‌ها را نمایه و فهرست می‌کنند. اگر یک URL به‌تازگی ساخته شده باشد، در کمتر از دو هفته در دستگاه نمایه‌سازی گوگل، نمایه می‌شود. وبگاه‌های نمایه شده را نیز هر ماه یکبار، با هدف درج تغییرات، از نو نمایه‌سازی می‌کند. این جریان نمایه‌سازی، بسیار سریع‌تر از موتورهای دیگر انجام می‌گیرد. (۳)

گوگل، روزانه به ۴۰ میلیون جستجو پاسخ می‌دهد (۱۱: ۱۰۸). سرعت بازیابی در موتور گوگل، کمتر از ۰/۵ ثانیه است. این موتور ادعا دارد که امکان جستجو به ۳۵ زبان مختلف را نیز فراهم کرده است. (۱)

نتایج جستجو در موتور گوگل از لحاظ میزان ارتباط به نسبت موتورهای دیگر تقریباً در بالاترین سطح قرار دارد. در تحقیقی که دیمتریوس و گوتلیب انجام دادند و موتور گوگل، یاهو و لایکوز را مقایسه کردند، در جستجو با کل واژگان، جستجوی مجاورتی و مکانی ساده و بولی، موتور گوگل از لحاظ میزان ارتباط، بالاترین امتیاز را کسب کرد. (۵: ۴۲)

موفقیت گوگل در بازیابی اطلاعات، از در اختیار داشتن بهترین الگوریتم جهت رتبه‌بندی نتایج بازیابی ناشی شده که بر ماهیت کاملاً آزادانه وب استوار است.

رتبه‌بندی صفحات وب در موتورهای جستجو

رتبه‌بندی صفحات بازیابی شده وب در پاسخ به هر جستجو، عامل مهمی است که می‌تواند

رضایت کاربران و در نهایت موفقیت موتور جستجو را تضمین کند. «با توجه به نتایج تحقیقات در این زمینه که نشان می‌دهد کاربران از چندین صفحه بازیابی کرده موتورهای جستجو، تنها تعداد محدودی از صفحات اولیه را مرور می‌کنند» (۴: ۳۹)، مسئله رتبه‌بندی نتایج، اهمیت دوچندان پیدا می‌کند.

سیاهه نتایج بازیابی شده، براساس معیارهای متنوع خاص هر موتور جستجو رتبه‌بندی می‌شود. از جمله این معیارها، می‌توان به تعداد واژگان منطبق شده با واژگان جعبه جستجو، مجاورت واژگان، محل قرار گرفتن واژگان در مدرک، بسامد واژگان (در هر مدرک و در کل پایگاه) و طول مدرک اشاره کرد (۴: ۳۹). کیم^۲ (۲۰۰۵) نیز، در مقاله خود موارد دیگری را یادآور می‌شود، مثل تعداد دفعاتی که هر رکورد را دیگر رکوردهای پایگاه مورد ارجاع قرار می‌دهند و همچنین، به نسبت اصطلاحات مرتبط با تعداد کل اصطلاحات موجود در رکورد اشاره کرده است. (۷: ۴۹)

فرمول دقیقی که طرز کاربرد این معیارها را نشان دهد، الگوریتم (خوارزمی) رتبه‌بندی نامیده می‌شود که از یک موتور جستجو به موتور جستجوی دیگر متفاوت است. به دلیل رقابت‌های بین‌المللی در میان شرکت‌های عرضه‌کننده موتورهای جستجو، برخی از این شرکت‌ها الگوریتم دقیق رتبه‌بندی را افشا نمی‌کنند (۹) و در عوض، الگوریتم‌ها و معیارهایی را که به کار می‌برند، به گونه‌ای

2. Kim.

کلی‌تر و فارغ از جزئیات در دسترس قرار می‌دهند.

موتور جستجوی هات بوت^۲، بسامد و جایگاه واژگان را به‌عنوان عامل‌های اولیه در نظر می‌گیرد. مدارکی که بسامد بالاتری از واژگان مورد جستجو دارند، وزن بیشتری به خود اختصاص می‌دهند. بسامد واژگان در کل پایگاه هم بر این مسئله تأثیرگذار است. علاوه بر این، تعداد تکرار ویژه نیز نسبت به طول متن در نظر گرفته می‌شود. زمانی که تعداد بسامد واژه برابر باشد، مدارک کوتاه نسبت به مدارک طولانی رتبه بالاتری کسب می‌کنند. همچنین وجود اصطلاحات مورد جستجو در عنوان مدرک و یا در ابربرچسب‌ها^۴، به نسبت زمانی که در متن مدرک وجود دارند، وزن بیشتری به مدارک اختصاص می‌دهند. (۴: ۳۹)

آلتا ویستا^۵ نیز، عامل‌های بیان شده را، به نحوی مورد ملاحظه قرار داده و علاوه بر این، تعداد اصطلاحات منطبق شده با واژگان مورد جستجو و همچنین مجاورت واژگان را در نظر داشته است. (۸: ۳۰)

پاره‌ای دیگر از موتورهای جستجو، اطلاعات کمتری در مورد نحوه رتبه‌بندی و عامل‌های آن، ارائه می‌کنند و فقط به تعدادی از عناصر اشاره دارند. از دیدگاه اینفوسیک^۶، وجود اصطلاحات مورد جستجو در عنوان و ابربرچسب ارزش فوق‌العاده‌ای

دارد. این، درحالی است که در نگاه لایکوز^۷، اصطلاحات موجود در عنوان و سرعنوان دیده نمی‌شود و اصطلاحات موجود در ابربرچسب، فاقد ارزش است (۴: ۳۹). اکسایت^۸ هم، اصطلاحات ابربرچسب را در نظر نمی‌گیرد و علاوه بر بازیابی مدارکی که اصطلاحات جستجو را شامل می‌شود، این موتور جستجو، محتوای مدارک را برای اصطلاحات مرتبط تحلیل می‌کند و این کار، طی فرایندی به نام ICE^۹ انجام می‌گیرد. (۴: ۳۹)

اخیراً پاره‌ای از موتورهای جستجو رویکردهای دیگری مانند توجه به لینک (پیوند)‌های رسیده و خارج شده از مدارک مدنظر قرار داده‌اند که از آن میان، موتور جستجوی گوگل را می‌توان نام برد.

طرح تحقیقاتی گوگل

در اوایل سال ۱۹۹۸، دانشجویان دانشگاه استنفورد کالیفرنیا، لارنس پیج و سرجی برین^{۱۰}، طی طرحی تحقیقاتی، الگوریتم (خوارزمی) خاصی برای رتبه‌بندی نتایج جستجوها براساس ارتباط بیشتر، طراحی کردند و آن را در مقاله‌ای با عنوان «کالبدشناسی موتور جستجوی وب فرامتنی در مقیاس بزرگ»^{۱۱} در همایش World Wide Web ارائه دادند.

این نوآوری، در ژانویه ۱۹۹۸ به نام لارنس پیج، به شماره ۶۲۸۵۹۹۹ به ثبت رسید و در

- HotBot.
- Meta tags.
- AltaVista.
- InfoSeek.
- Lycos.

- Excite.
- Intelligent Concept Extraction.
- Lawrence Page & Sergey Brin.
- The anatomy of a large-scale hypertextual; web search engine.

سپتامبر ۲۰۰۱ مورد حمایت مالی قرار گرفت و بدین ترتیب، شرکت گوگل به طور رسمی پایه‌گذاری و تأسیس شد. الگوریتم معرفی‌شده، با عنوان PageRank شناخته شد و روزی ۱۵۰ میلیون درخواست را پاسخ می‌گفت. امروزه گوگل با این الگوریتم (خوارزمی)، موفقیت خود را تضمین کرده است. (۵: ۴۲)

الگوریتم (خوارزمی) بازیابی و رتبه‌بندی گوگل

در الگوریتم گوگل، دو بخش اصلی وجود دارد.

بخش اول، سامانه text matching نامیده می‌شود که در آن، موتور گوگل صفحات مرتبط با واژه‌های وارد شده در جعبه جستجو را بازیابی می‌کند.

بخش دوم - که از اهمیت معادلی برخوردار است - سامانه PageRank نامیده می‌شود که طی آن، صفحات ثبت شده گوگل رتبه‌بندی می‌گردد. (۳)

سامانه text matching

گوگل در جستجوی اطلاعات مربوط، ارزش و وزن زیادی به برچسب عنوان می‌دهد. این که واژه‌ها یا اصطلاحات مورد جستجو در برچسب عنوان وجود داشته باشد، امر درخور توجهی است. «گوگل، برچسب‌های شبیه کلیدواژه‌ها را در نظر نمی‌گیرد. این، بدان دلیل است که گوگل احتمال می‌دهد که مدیران وبگاه‌ها از این عامل سوءاستفاده کنند و به منظور جلب ملاقات‌کنندگان بیشتر،

شماری واژه‌های نامربوط در این بخش وارد نمایند» (۲). در مقابل، گوگل، توصیفات و واژگان را از چند خط اول متن روی صفحه وبگاه اخذ می‌کند؛ یعنی اگر کلیدواژه‌های مرتبط در بالای صفحه وب باشند، آن صفحه به عنوان مربوط شناخته می‌شود. (۳)

گوگل، تجمع واژگان مورد جستجو در بدنه اصلی متن و صفحه را نیز، به عنوان نشانه ارتباط در نظر می‌گیرد. تجمع ۶ تا ۱۰ درصدی، بهترین حالت ممکن است. (۳)

گوگل، اخیراً توجه خاصی به برچسب‌های سرعنوان‌های متعدد ظاهر شده در بخش‌های متن کرده است و به آن‌ها ارزش داده است. ابزار دیگر گوگل برای انتخاب صفحه‌های ذی‌ربط، در نظر گرفتن کلیدواژه‌های مربوط در برچسب‌های برجسته و مشخص است. (۳)

سامانه PageRank

مفهوم PageRank؛ از مراحل آغازین وب جهان گستر، هر موتور جستجو، شیوه‌های متنوعی را جهت رتبه‌بندی نتایج جستجو توسعه داده است. در همین راستا و به منظور دستیابی به اهداف رتبه‌بندی مبتنی بر محتوا، موتور گوگل رتبه‌بندی مبتنی بر میزان لینک‌ها را گسترش داد.

برای این منظور، تعداد لینک‌هایی که از متون دیگر به متن حاضر وارد می‌شود، به عنوان نقطه اهمیت برای متن و صفحه حاضر به‌شمار می‌آید. هر لینک (پیوند) به صفحه حاضر، به عنوان یک رأی مثبت برای آن شمرده می‌شود و هر صفحه وبی - که لینک‌های بیشتری به آن وارد می‌گردد - با

اهمیت‌تر تلقی می‌گردد.

برخلاف مفهومی که از سامانه لینک‌ها استنباط می‌گردد، محاسبه ارزش هر صفحه وب به همین سادگی هم نیست و متغیرهای دیگری هم در آن دخالت دارد. هر صفحه‌ای که به صفحه حاضر لینک (پیوند) داشته باشد، به یک اندازه به صفحه حاضر ارزش نمی‌دهد، بلکه ارزش هر لینک برای سند زمانی مهم است که آن لینک از صفحه باارزشی باشد و آن صفحه در الگوی رتبه‌بندی صفحات، از رتبه بالایی برخوردار باشد. بنابراین، رتبه هر سند به رتبه دیگر اسنادی بستگی دارد که به آن لینک می‌دهند و رتبه آن اسناد هم به همین ترتیب تعیین می‌گردد. این چرخه همچنان ادامه می‌یابد تا اینکه رتبه هر سند به رتبه کل وب بستگی پیدا کند و این، همان ماهیت آزادانه این الگوست. (۱۰)

متغیر دیگری که در این الگوریتم (خوارزمی) دخالت دارد، میزان لینک‌های دیگری است که از صفحه‌ای که به صفحه حاضر لینک دارد، خارج می‌گردد. مثلاً اگر صفحه ما - که قرار است رتبه‌اش محاسبه شود - A فرض گردد و صفحه‌ای که به همین صفحه A لینک داده است، B فرض شود، چنانچه تنها لینکی که از B خارج گردیده، لینک به A باشد، این لینک ارزش بالاتری نسبت به موقعی دارد که از B، ۲۰ لینک خارج گردیده باشد و لینک به A هم یکی از آن ۲۰ لینک بوده باشد.

یعنی هر چقدر صفحه‌ای که به صفحه ما لینک داده یا به اصطلاح رأی داده،

لینک‌های بیشتری داشته باشد، رتبه کمتری (PageRank) به صفحه ما تعلق می‌گیرد و برعکس. (۶)

اگرچه این رهیافت به نظر خیلی گسترده و پیچیده می‌آید، اما پیچ و برین، موفق شدند این را با استفاده از الگوریتم نوآورانه خود عملی سازند.

الگوریتم (خوارزمی) PageRank

فرمول الگوریتم گوگل - که آن را پیچ و برین طراحی کردند - به این صورت است:

$$PR(A) = (1-d) + d \left(\frac{PR(Ti)}{C(Ti)} + \dots + \frac{PR(Tn)}{C(Tn)} \right)$$

PR (A) : PR یا رتبه صفحه A است.

PR (Ti) : PR صفحه Ti که به صفحه A لینک داده است.

C (Ti) : تعداد لینک‌هایی که از صفحه

Ti خارج می‌شود.

d : یک damping factor که بین صفر و یک است.

از این فرمول، رتبه‌بندی موارد زیر دریافت می‌گردد:

- این سامانه، صفحات وب را به‌طور کلی رتبه‌بندی نمی‌کند، بلکه رتبه هر صفحه به تنهایی محاسبه می‌گردد و رتبه هر صفحه، به وسیله دیگر صفحاتی که به آن صفحه لینک داده‌اند، محاسبه می‌گردد.

- هرچه تعداد لینک‌های وارد شده به صفحه حاضر بیشتر باشد، رتبه صفحه حاضر بالاتر است.

- هرچه رتبه صفحه لینک داده به صفحه

۱۲. کاربرد PR در مقاله، معادل PageRank است.

حاضر بیشتر باشد، رتبه صفحه حاضر بالاتر می‌رود.

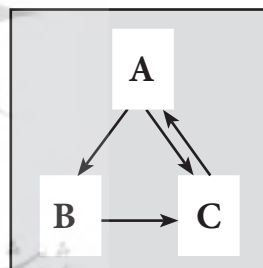
- هرچه تعداد کل لینک‌های صفحه لینک داده به صفحه حاضر بیشتر باشد، رتبه کمتری به صفحه حاضر می‌دهد.

یعنی یک وبگاه با رتبه (PR) ۲ و تعداد ۲ لینک بیش از وبگاهی با رتبه ۶ و تعداد ۱۰ لینک به وبگاه ما ارزش می‌دهد.

مقیاس PageRank (PR)، از یک تا ده است. وبگاه کم اهمیت PR برابر با یک تا سه دارد و وبگاه پراهمیت PR، از هفت تا ده می‌تواند داشته باشد. (۳)

مثال:

یک وب کوچک را، مشتمل بر سه صفحه A، B و C در نظر می‌گیریم:



صفحه A لینک به هر دو صفحه B و C، صفحه B، لینک به صفحه C، و صفحه C، لینک به صفحه A دارد. برطبق نظر پیج و برین، damping factor معمولاً ۰/۸۵ است. اما برای ساده‌تر شدن محاسبه، در اینجا آن را ۰/۵ در نظر می‌گیریم. میزان دقیق damping factor بر PR تأثیرگذار است، ولی بر اصول اساسی آن تأثیر ندارد.

براساس فرمول کلی زیر:

$$PR(A) = (1-d) + d(PR(Ti) / C(T) + \dots + PR(Tn) / C(Tn))$$

رتبه سه صفحه را محاسبه می‌کنیم:

$$PR(A) = 0.5 + 0.5 PR(C)$$

$$PR(B) = 0.5 + 0.5 (PR(A) / 2)$$

$$PR(C) = 0.5 + 0.5 (PR(A) / 2 + PR(B))$$

بنابراین:

$$PR(A) = 14/13 = 1.07692308$$

$$PR(B) = 10/13 = 0.76923077$$

$$PR(C) = 15/13 = 1.15384615$$

روشن است که جمع PR تمام صفحات، ۳ می‌شود؛ بنابراین معادل با مجموع کل صفحات وب است^{۱۳}. تعیین کردن این ارزش برای این سه صفحه راحت بود، ولی در نظر بگیریم که این عملیات در کل وب - که شامل میلیون‌ها سند است - چگونه انجام می‌گیرد؟

موتور جستجوی گوگل، با به خدمت گرفتن شاخه‌های مختلف علم ریاضی از جمله جبر خطی، انجام این عملیات و محاسبات را به بهترین نحو ممکن عملی کرده است.

دیگر عامل‌های تأثیرگذار بر PageRank

از زمان کار علمی سرجی برین و لورنس پیج درباره PageRank، مباحث زیادی در زمینه عامل‌های دیگری که علاوه بر دستورالعمل‌های لینک‌ها بر PR تأثیرگذار می‌باشد، به وجود آمده است. خود لورنس پیج، عامل‌های تأثیرگذار بالقوه‌ای به شرح زیر بیان کرده است:

۱۳. برای اطلاعات بیشتر در مورد مثال، به منبع شماره ۱۰ بنگرید.

۱. لینک‌های قابل مشاهده؛

لینک‌های قابل مشاهده، شامل لینک‌هایی است که برجسته‌ترند و توجه را جلب می‌کنند که از آن میان، به لینک‌های رنگی و یا نوشته شده با قلم ایتالیک، می‌توان اشاره کرد. این لینک‌ها به دلیل اینکه احتمال تلیک کردن تصادفی کاربر روی آن‌ها بیشتر است، از ارزش بیشتری برخوردارند و رتبه بالاتری برای صفحه‌ای که لینک به آن می‌رسد، به وجود می‌آورند.

۲. موقعیت لینک‌ها؛

عامل دیگری که می‌تواند بر ارزش لینک تأثیرگذار باشد، موقعیت قرار گرفتن لینک در صفحه است. لینک‌هایی که در نیمه بالای صفحه قرار دارد، نسبت به لینک‌هایی که در نیمه پایین صفحه قرار دارند، از ارزش بیشتری برخوردارند و برعکس. لینک‌های با ارزش‌تر، برای صفحه مورد لینک ارزش بیشتری به وجود می‌آورند.

۳. فاصله بین صفحات وب؛

مدیران صفحات وب، ممکن است تأثیرات غیرواقعی در عملکرد PR ایجاد کنند و آن، بدین صورت است که مدیران وبگاه‌ها می‌توانند تعداد کثیری از صفحات وب را ایجاد کنند و لینک‌های متقابلی در بین آن‌ها توزیع نمایند که موجب می‌شود رتبه صفحه‌ای خاص به‌طور غیرواقعی در الگوریتم بالا برود. در این حالت، صفحات وب بدون اینکه از صفحات با PR بالا لینک

داشته باشند، می‌توانند رتبه بالایی به خود اختصاص دهند^{۱۴}. در این حالت، نه فقط مفهوم PR تضعیف می‌گردد، بلکه نمایه موتور جستجو نیز با میزان بی‌شماری از صفحات وب بدون کیفیت که فقط برای تأثیرگذاری بر PR به وجود آمده‌اند، اشباع می‌گردد. از این رو، لورنس پیج در جزئیاتی که در سندش ارائه داده، ارزش خاصی به لینک‌هایی می‌دهد که از صفحاتی با مسافت دورتر می‌رسد و این امر، می‌تواند از تأثیرات غیرواقعی بر PR جلوگیری کند. چون با وجود مسافت بین صفحات، مدیران وبگاه‌ها نمی‌توانند نظارت زیادی بر آن‌ها داشته باشند، یک مقیاس برای مسافت بین دو صفحه می‌تواند براساس Domain باشد. در این حالت، لینک‌های رسیده از داخل Domain، ارزش کمتری نسبت به لینک‌های خارج از آن دارند و لینک‌هایی با ارزش‌ترند که از صفحه‌ای خارج از Domain مربوط باشند. برای مسافت، مقیاس‌های دیگری هم می‌توان در نظر گرفت که از آن میان، می‌توان به فاصله جغرافیایی میان خدمتگرها اشاره کرد.

۴. روز آمد بودن صفحات لینک‌دهنده؛

لورنس پیج، در این باره نیز بیان داشته که روزآمد بودن صفحات و اسنادی که لینک را ایجاد کرده‌اند، می‌تواند مبین این باشد که اطلاعات مندرج در صفحه - که لینک به آن می‌رسد - تا چه حد روزآمدند و ارزش بیشتری برای صفحه ایجاد می‌کنند.

۱۴. شاید بتوان این فرایند را با self citation در نوشته‌های چاپی معادل دانست. مؤسسه اطلاعات علمی (ISI) نیز، برای رتبه‌بندی مجلات براساس میزان استناد به مقالات آن‌ها با همین مشکل خود استنادی روبه‌روست و در فرمول جداگانه، موارد خوداستنادی را حذف می‌کند.

الگوی موج سوار تصادفی: تأییدی بر الگوریتم (خوارزمی)

لورنس پیج و سرچی برین، تأیید ساده و مستقیمی برای الگوریتمشان ارائه کرده‌اند. آن‌ها PR را به عنوان الگویی از رفتار کاربری در نظر گرفته‌اند که با موج سواری در وب، روی لینک‌ها بدون توجه به محتوای آن‌ها تلیک می‌کند. کاربر با موج سواری، صفحه خاصی را با یک احتمال مطمئن ملاقات می‌کند که این، از PR صفحه سرچشمه می‌گیرد. احتمال اینکه او به طور تصادفی روی یک لینک در یک صفحه تلیک کند، به مدد تعداد لینک‌های روی آن صفحه تعیین می‌گردد. به همین دلیل است که در فرمول PR، هر صفحه به تعداد لینک‌های روی صفحه تقسیم می‌گردد. بنابراین، احتمال اینکه کاربر به طور تصادفی به یک صفحه برسد، برابر با جمع احتمالاتی است که در برابر او قرار می‌گیرد تا روی لینک‌ها تلیک کند و به یک صفحه خاص برسد. البته اکنون این احتمال به وسیله d کاهش یافته است. تأیید بر الگوی موج سوار تصادفی^{۱۵}، این است که موج سوار در وب روی تعداد مشخصی از لینک‌ها تلیک نمی‌کند، بلکه گاهی خسته می‌شود و به طور تصادفی به صفحه دیگری پرش می‌کند.

احتمال برای موج سوار تصادفی با تلیک بر لینک‌های داده شده، به وسیله d damping factor - که بسته به درجه احتمال، بین صفر و یک متغیر است - متوقف نمی‌شود. زمانی که فرد روی لینک‌های تلیک شده باقی

می‌ماند، d بالاتر است و زمانی که تصادفی به صفحه دیگر پرش می‌کند، بعد از اتمام لینک‌های تلیک شده، احتمال به عنوان یک $(1-d)$ دائم در الگوریتم باقی می‌ماند. همچنین با توجه به لینک‌های وارد شده، احتمال پرش تصادفی کاربر به یک صفحه، همیشه $(1-d)$ است. بنابراین، هر صفحه، همواره یک حداقل PageRank را داراست. (۱۰)

ویرایش دوم الگوریتم (خوارزمی) گوگل

دومین ویرایش الگوریتم - که به همت لورنس پیج و سرچی برین منتشر شده - به این قرار است:

$$PR(A) = (1-d)/N + d(PR(T_1)/C(T_1) + \dots + PR(T_n)/C(T_n))$$

که در اینجا، N مجموع کامل صفحات وب است. دومین ویرایش، تفاوت بنیادی با ویرایش اول ندارد. در این ویرایش براساس الگوی موج سوار تصادفی، PageRank یا رتبه‌بندی صفحه، یک احتمال واقعی برای موج سوار در وب است که براساس آن، او بعد از تلیک کردن روی لینک‌های متعدد به صفحه‌ای خاص می‌رسد. کل PageRank، یک توزیع احتمال را در کل صفحات وب شکل می‌دهد. بنابراین، جمع کل PageRank صفحات وب، برابر با یک خواهد شد.

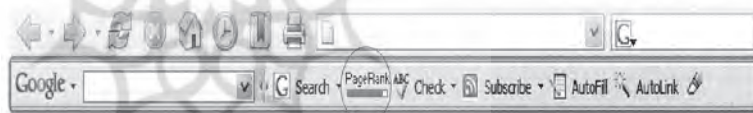
برعکس، در اولین ویرایش الگوریتم، احتمال برای رسیدن تصادفی کاربر به یک صفحه خاص، به کمک جمع تعداد صفحات وب ارزش‌گذاری می‌شود. بنابراین در این

15. The random surfer model.

ویرایش، PR ارزش درخور انتظاری برای موج‌سواری تصادفی کاربری خواهد بود که صفحه خاصی را ملاقات می‌کند و این روند را به هر تعداد صفحه‌ای که وب داشته باشد، دوباره از نو شروع می‌کند. مثلاً اگر وب یکصد صفحه داشته باشد، و یک صفحه با PR برابر با ۲ داشته باشیم، موج‌سواری تصادفی، کاربر را به‌طور متوسط ۲ بار به آن صفحه خواهد رساند، حتی اگر او یکصد بار از نو شروع کند.

نمایش PageRank از نوار ابزار گوگل

نوار ابزار گوگل، به‌عنوان مرورگری است که می‌تواند از وبگاه گوگل دریافت گردد و تعدادی از ویژگی‌ها را برای جستجوی بهتر و راحت‌تر فراهم کند.



این نوار ابزار، PageRank را روی مقیاس ۱ تا ۱۰ نمایش می‌دهد و این، یکی از خصوصیات پیشرفته گوگل است. یعنی کاربر در هر صفحه‌ای از وب که حضور داشته باشد، و از طریق موتور جستجوی گوگل آن صفحه خاص را باز کرده باشد، می‌تواند با تلیک کردن بر PageRank - که بر روی نوار ابزار قرار دارد - رتبه صفحه جاری را مشاهده کند.^{۱۶}

نتیجه‌گیری

گوگل، بهترین بودن را به‌عنوان نقطه پایان قبول ندارد، بلکه آن را نقطه آغاز می‌داند. اکنون گوگل بر آن است تا دست کم در زمینه موتورهای جستجو پیشرو باشد و برای رسیدن به این هدف، پیوسته نوآوری می‌کند و فن‌آوری خود را گسترش می‌دهد. نمونه بارز آن، الگوریتم (خوارزمی) رتبه‌بندی گوگل است که می‌توانیم بگوییم در عمل بی‌نظیر است و گوگل نیز به آن می‌بالد. این الگوریتم، رتبه‌بندی صفحات را برپایه میزان و کیفیت لینک‌ها قرار داده است و کمیت و کیفیت را بر پایه منطق ریاضی محاسبه می‌کند. او این رتبه‌بندی را بر پایه دموکراسی وب می‌داند و تا حدود زیادی هم می‌توانیم بگوییم نسبت به موتور جستجوهای دیگر موفقیت‌های زیادی

کسب کرده است. ولی نکته شایان ذکر در اینجا، این است که با وجود موقعیت کارآ و درخشان این ابزار از لحاظ فنی و عملی، نباید از نظر دور داشت که هنوز هم برای برخی از موضوعات، دقت مضمونی از الگوریتم ریاضی بااهمیت‌تر است.

منابع

۱. پژوهشگاه اطلاعات و مدارک علمی ایران.

[قابل دسترسی در:]

۱۶. جهت دان لود مرورگر نوار ابزار گوگل، در صفحه خانگی گوگل، روی گزینه more تلیک می‌کنیم، و بعد از سیاه‌های که باز می‌شود، گزینه toolbar را تلیک کنیم و بعد، از منوی بعدی گزینه download را برگزینیم.

7. Kim, Guenther. "Choosing a search engine". *Online*, Vol.29, No.1 (Jan./Feb 2005): 48.

8. Ressel, M Robert. "Choosing a search engine for internet exploration". *The Technology Teacher*, No.56 (Apr.1997): 30.

9. Sulivon, Danny. "How search engines rank web pages". [on-line]. Available: http://search_engine_watch.com/web_masters/article.php/2167967.

10. "A survey of google pageRank". [on-line]. Available: <http://www.miswebdesign.com/resources/articles/pageRank-8.html>.

11. Timothy, Archibald. "Search us, says google". *Technology Review*, No.3 (Nov./Dec.2000): 108.

<http://www.irandoc.ac.ir/ETELAART/16/16-3-4-7-htm>.

2. Brandt, Daniel. "PageRank: google original sin". [on-line]. Available: <http://www.google-watch.org/pagerank.html>.

3. Callan, David. "Google ranking tips". [on-line]. Available: <http://www.akamarketing.com/google-ranking>.

4. Courtois, Martin P; Berry, Michael W. "Result ranking in web search engines". *Online*, Vol.23, No.3 (May/Jun.1999): 39.

5. Demetios, Eliopoulos; Gotlieb, Calvin. "Evaluating web search ranking". *Online*, Vol.27, No.2 (Mar./Apr. 2003): 42.

6. Google technology. [on-line]. Available: <http://www.google.com/technology/index.html>

تاریخ دریافت: ۱۳۸۵/۵/۸

پژوهشگاه علوم انسانی و مطالعات فرهنگی

